

## PATENT ABSTRACTS OF JAPAN

(11)Publication number : 08-077010

(43)Date of publication of application : 22.03.1996

(51)Int.Cl. G06F 9/44  
G06F 17/60  
G06F 19/00

(21)Application number : 06-239437

(71)Applicant : HITACHI LTD

(22)Date of filing : 07.09.1994

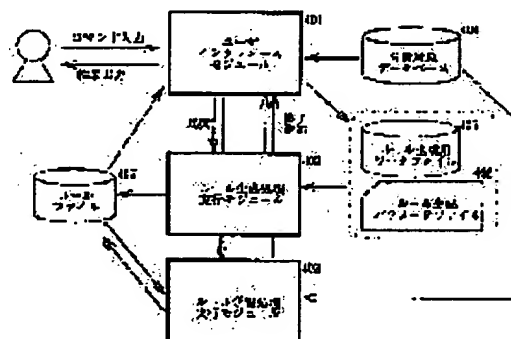
(72)Inventor : MAEDA AKIRA  
ASHIDA HITOSHI  
TANIGUCHI YOJI  
ITO YUKIYASU  
TAKAHASHI YORI

## (54) METHOD AND DEVICE FOR DATA ANALYSIS

## (57)Abstract:

**PURPOSE:** To provide excellent operation environment wherein a user can generate a result that precisely represent features of data without having any special knowledge when extracting regularity and cause-effect relation between data effective for the user in the form of a rule from information stored in a data base.

**CONSTITUTION:** The user specifies and corrects a parameter required for processing through a user interface module 401 and utilizes a displayed rule generation result. A rule generating process execution module 402 outputs a rule file 407 based upon input data 405 by using the process parameter which is specified or automatically determined according to data. A rule learning process execution module 403 optimizes the rule by using the contents of a data base 404 to be analyzed and a rule file 407 and stores the result in the rule file 407.



## LEGAL STATUS

[Date of request for examination] 23.07.1997

[Date of sending the examiner's decision of rejection] 24.03.2000

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2000 Japanese Patent Office

(19) 日本国特許庁 (JP)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平 8 - 7 7 0 1 0

(43) 公開日 平成8年(1996)3月22日

(51) Int. Cl.<sup>6</sup>

G 0 6 F 9/44  
17/60  
19/00

識別記号 庁内整理番号  
5 5 0 N 7737 - 5 B

F I

技術表示箇所

G 0 6 F 15/21  
15/30

Z  
Z

審査請求 未請求 請求項の数 1 3

F D

(全 2 0 頁)

(21) 出願番号 特願平6-239437

(22) 出願日 平成6年(1994)9月7日

(71) 出願人 000005108

株式会社日立製作所

東京都千代田区神田駿河台四丁目6番地

(72) 発明者 前田 章

神奈川県川崎市麻生区王禅寺1099番地 株

式会社日立製作所システム開発研究所内

(72) 発明者 芦田 仁史

神奈川県川崎市麻生区王禅寺1099番地 株

式会社日立製作所システム開発研究所内

(72) 発明者 谷口 洋司

神奈川県川崎市麻生区王禅寺1099番地 株

式会社日立製作所システム開発研究所内

(74) 代理人 弁理士 矢島 保夫

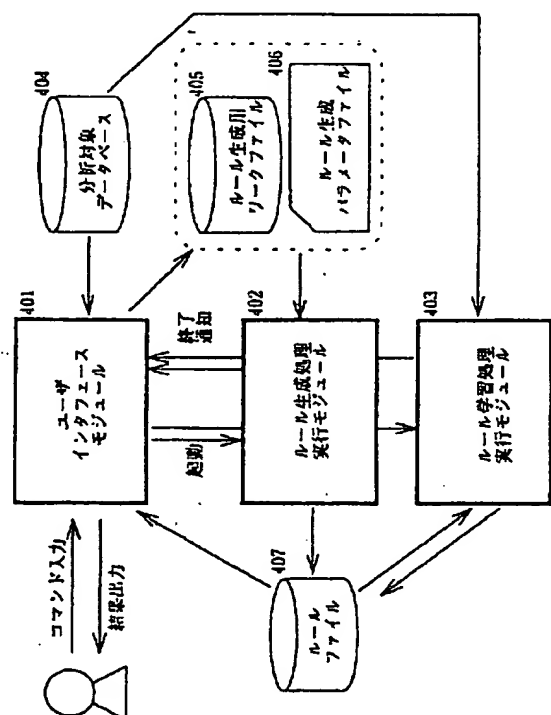
最終頁に続く

(54) 【発明の名称】 データ分析方法および装置

(57) 【要約】 (修正有)

【目的】 データベースに蓄積された情報から、利用者にとって有効なデータ間の規則性・因果関係をルール形式で抽出するのに、利用者が特別な知識を持たなくてもデータの特徴を精度よく表現する結果を生成できる良好な操作環境を提供する。

【構成】 利用者は、ユーザインタフェースモジュール 401 を介して処理に必要なパラメータの指定・修正を行い、また表示されたルール生成結果を利用する。ルール生成処理実行モジュール 402 は、指定された、またはデータから自動的に決定した処理パラメータを用いて、入力データ 405 からルールファイル 407 を出力する。ルール学習処理実行モジュール 403 は、分析対象データベース 404 とルールファイル 407 の内容を用いて、ルールの最適化を行い、結果をルールファイル 407 に格納する。



## 【特許請求の範囲】

【請求項 1】複数のデータ項目の具体的な値である数値または記号値のデータの集まりを 1 つのレコードとし、複数のレコードからなるデータを入力して分析することにより、該データの特徴を表現するルールを生成出力するデータ分析方法であって、前記複数のデータ項目中から、ルールの条件節と結論節に用いるデータ項目を選択するステップと、

選択したデータ項目が数値の値を持つデータ項目であるときは、数値を記号値に変換するステップと、データ項目間の相互関係を項目名称と記号値の組みを 1 つまたは複数個組み合わせたルール形式で表現した複数の候補ルールを作成するステップと、該候補ルールが表現する相互関係の強さを評価する評価尺度を算出するステップと、作成した複数の候補ルールのうち、該評価尺度が大きな候補ルールを 1 つまたは複数個求めるステップと、求めた候補ルールを出力するステップとを備えたことを特徴とするデータ分析方法。

【請求項 2】複数のデータ項目の具体的な値である数値または記号値のデータの集まりを 1 つのレコードとし、複数のレコードからなるデータを入力して分析することにより、該データの特徴を表現するルールを生成出力するデータ分析方法であって、前記複数のデータ項目中から、ルールの結論節に用いるデータ項目である結論項目を 1 つ選択するステップと、

前記複数のデータ項目中から、ルールの条件節に用いるデータ項目である条件項目を選択するステップと、

選択したデータ項目が数値の値を持つデータ項目であるときは、数値を記号値に変換するステップと、

条件項目の項目名称とその記号値の組みを 1 つまたは複数個組み合わせる条件節を作成し、結論項目の項目名称とその記号値を組み合わせる結論節を作成し、作成した条件節と結論節を組み合わせる複数の候補ルールを作成するステップと、

作成した候補ルールの良さを評価する評価尺度を算出するステップと、

作成した複数の候補ルールのうち、算出した評価尺度が大きな候補ルールを 1 つまたは複数個求めるステップと、

求めた候補ルールを最終的に生成したルールとして出力するステップとを備えたことを特徴とするデータ分析方法。

【請求項 3】前記ルールの条件節に用いるデータ項目を選択するステップは、各データ項目間の関連度を求め、その関連度に基づいて、前記ルールの条件節に用いるデータ項目の選択を自動的に行う請求項 1 または 2 のいずれか 1 つに記載のデータ分析方法。

【請求項 4】前記条件項目を選択するステップは、前記複数のデータ項目中の前記結論項目以外のデータ項目

目のすべての組み合わせに対して、それぞれ、項目間の関連性の尺度を示す類似度尺度を計算するステップと、計算した類似度尺度にしたがって、関連性が高い前記結論項目以外のデータ項目を集めて、幾つかのクラスタにクラスタリングするステップと、

前記結論項目と前記結論項目以外のデータ項目とのすべての組み合わせに対して、それぞれ、項目間の関連性の尺度を示す依存度尺度を計算するステップと、

計算した依存度尺度にしたがって、前記クラスタ中で前記結論項目と最も関連性が高いデータ項目を選び、前記条件項目として採用するステップとを備えた請求項 2 に記載のデータ分析方法。

【請求項 5】前記数値を記号値に変換するステップは、幾つかの記号値と前記数値が該記号値に適合する度合を示す記号値ごとのメンバシップ関数とを決定することにより行う請求項 1 または 2 のいずれか 1 つに記載のデータ分析方法。

【請求項 6】前記メンバシップ関数値に基づいて、入力されたレコードの重みを算出し、該重みの情報を用いてルールの相互関係の強さを評価する請求項 5 に記載のデータ分析方法。

【請求項 7】分析すべき前記レコード中のあるデータ項目の値が欠損値を含むとき、該データ項目の記号値として欠損値を示す記号値を追加し、生成するルールに該欠損値を示す記号値を使用するか否かを各項目ごとに指定するステップを、さらに備えた請求項 1 または 2 のいずれか 1 つに記載のデータ分析方法。

【請求項 8】前記欠損値を示す記号値を使用するか否かの指定を、そのデータ項目と結論項目との依存関係に基づいて自動的に決定するステップを、さらに備えた請求項 7 に記載のデータ分析方法。

【請求項 9】前記候補ルールの評価尺度は、該候補ルールが対象とするレコード数で表現されるカバー率と、該候補ルールの正解率で表現される精度とから算出され、かつカバー率と精度を重要視する度合をパラメータとして指定するステップをさらに備えた請求項 1 または 2 のいずれか 1 つに記載のデータ分析方法。

【請求項 10】入力レコードの前記データ項目間に存在する関数関係を求めるステップをさらに設け、該関数関係に基づいて冗長なルールの生成および表示を行わない請求項 1 または 2 のいずれか 1 つに記載のデータ分析方法。

【請求項 11】出力されたルールの全部または一部と、入力されたデータを用いて、前記ルールの評価尺度を最大化するように、前記メンバシップ関数の形状を指定するパラメータの値を調整するステップを、さらに備えた請求項 5 または 6 のいずれか 1 つに記載のデータ分析方法。

【請求項 12】複数のデータ項目の具体的な値である数値または記号値のデータの集まりを 1 つのレコードと

し、複数のレコードからなるデータを入力して分析することにより、該データの特徴を表現するルールを生成出力するデータ分析装置であって、前記複数のデータ項目中から、ルールの条件節と結論節に用いるデータ項目を選択する手段と、

選択したデータ項目が数値の値を持つデータ項目であるときは、数値を記号値に変換する手段と、

データ項目間の相互関係を項目名称と記号値の組みを1つまたは複数個組み合わせたルール形式で表現した複数の候補ルールを作成する手段と、

該候補ルールが表現する相互関係の強さを評価する評価尺度を算出する手段と、

作成した複数の候補ルールのうち、該評価尺度が大きな候補ルールを1つまたは複数個求める手段と、

求めた候補ルールを出力する手段とを備えたことを特徴とするデータ分析装置。

【請求項13】複数のデータ項目の具体的な値である数値または記号値のデータの集まりを1つのレコードとし、複数のレコードからなるデータを入力して分析することにより、該データの特徴を表現するルールを生成出力するデータ分析装置であって、前記複数のデータ項目中から、ルールの結論節に用いるデータ項目である結論項目を1つ選択する手段と、

前記複数のデータ項目中から、ルールの条件節に用いるデータ項目である条件項目を選択する手段と、

選択したデータ項目が数値の値を持つデータ項目であるときは、数値を記号値に変換する手段と、

条件項目の項目名称とその記号値の組みを1つまたは複数個組み合わせ条件節を作成し、結論項目の項目名称とその記号値を組み合わせて結論節を作成し、作成した条件節と結論節を組み合わせる複数の候補ルールを作成する手段と、

作成した候補ルールの良さを評価する評価尺度を算出する手段と、

作成した複数の候補ルールのうち、算出した評価尺度が大きな候補ルールを1つまたは複数個求める手段と、

求めた候補ルールを最終的に生成したルールとして出力する手段とを備えたことを特徴とするデータ分析装置。

【発明の詳細な説明】

【0001】

【産業上の利用分野】本発明は、データベース装置などの情報記憶装置内に格納された数値または記号で表現されたデータの集まりを分析し、利用者にとって有用な表現に加工変換するデータ分析方法および装置に関する。

【0002】

【従来の技術】計算機技術の進歩により、計算機内に蓄積されるデータ量は年々増大している。特に、ネットワーク化が進むにつれて、オンラインシステムを中心にこの傾向はますます顕著になっている。現在では、レコード数で100万件、データ量でギガ(=10の9乗)バ

イトを越えるものも珍しくない。

【0003】計算機内に蓄積されたデータは、それだけでは数値や記号の集まりに過ぎないことが多い。そこで、このデータの集まりを利用者にとって有用な情報に変換し、データの有効活用を図ろうとする技術が提案されている。最も広く知られている手法は、相関分析・重回帰分析といった統計的手法である。

【0004】また、比較的新しい手法としては、利用者にとって理解しやすい「もし～ならば…である」といったルール形式に変換して結果を出力する方法、すなわちルールインダクションと呼ばれる知識獲得手法を用いる方法が、知られている。例えば、日立クリエイティブワークステーション2050(商品名)ES/TOOL/W-R I解説/操作マニュアルのページ23～31には、データ間に存在する関係をルールの形で表現する手法が述べられている。

【0005】この方法は、本来は与えられたデータからエキスパートシステムへ入力可能なルールを生成することを目的としたものであった。しかし、蓄積されたデータに含まれている因果関係や規則性といった特徴を、人間である利用者が発見するという目的に利用することが可能である。

【0006】

【発明が解決しようとする課題】上記従来技術は、あくまで計算機が利用できるルールを生成することを目的としている。したがって、人間である利用者が、そのルールを解釈することは可能ではあるが、もとより人間に理解しやすい形式でルールを生成する訳ではないので、人間が解釈しデータの特徴を理解するのに有効なルールを生成することはできなかった。例を用いて、上記文献に述べられている方法を説明する。

【0007】まず、データを個々の事例の集合と考える。例えば、半導体製造プロセスにおける品質管理データベースを用いて、半導体の不良要因の分析をするといった利用方法では、個々の事例はウェハとよばれる製造単位で管理され、それぞれの製造工程における処理パラメータや、各種検査結果などの情報の組を一つの事例として扱うことができる。図1に、このようなデータの例を示す。

【0008】また、銀行の顧客データベースから顧客の金融商品の購入動向を調べるといった利用方法では、年齢・預金残高・職業・年収・金融商品購入履歴といった顧客ごとの情報の組が1つの事例であり、分析の対象とするデータはこの事例の集まりと見なすことができる。この例に関しては、実施例で詳しく説明する。

【0009】上記従来技術によるルール生成の例を説明する。例として、ある金融商品(商品Aとする)を購入した顧客に共通の特徴を調べることを考える。この場合、従来技術では、データの各項目(年齢や預金残高など)の値から、金融商品Aを購入した顧客に対応する事

例と購入していない顧客に対応する事例とをできるだけ精度よく分類するルールを生成することが目的となる。

【0010】上記従来技術では、項目の値の組（例えば年齢が40歳以上でかつ預金残高が1000万以上、など）の中で、与えられたデータをもっとも精度よく分類する組を生成する。この場合の精度とは、特定の値を持つ事例の部分集合の中で、金融商品Aを購入した顧客に対応する事例の割合が大きければ大きいほど、購入した顧客の特徴を精度よく分類するものと考ええる。この値の組は、「IF 年齢が40歳以上 AND 預金残高が1000万以上 THEN 金融商品Aを購入」というルールの形式で表現することができる。

【0011】次に、生成されたルールで説明される事例を、全体の事例の集合から除く。上記の例では、年齢が40歳以上で預金残高が1000万以上という条件を満たす事例を除くことになる。そして、残された事例の集合に対して、もっとも精度よく分類する項目の値の組を求める。以上の処理を繰り返すことにより、金融商品Aを購入した顧客を、購入していない顧客から識別するためのルール群を得ることができる。

【0012】以上の説明から分かるように、上記従来技術で得られるルール群は、  
IF 年齢が40歳以上 AND 預金残高が1000万以上 THEN 金融商品Bを購入  
ELSE IF 職業が自営業 AND 年収が800万以上 THEN 金融商品Bを購入  
ELSE IF …というように、IF … ELSE IF … ELSE IF …の形式となる。

【0013】計算機がこのルール群を利用して分類をする場合には、先頭のIFから順に調べていくだけで機械的に処理を実行できる。しかしながら、ルールの数が増えれば増えるほど、人間がこのルールから金融商品Aを購入した顧客の特徴を理解することは困難になっていく。さらに、1つルールを生成する毎に、残された事例の集合からルールを探しに行くという処理を繰り返すため、事例の数が大きくなるにつれ、必要な処理時間が急速に増大することも問題である。

【0014】さらに重大な問題として、上記の例題のような実世界のデータの場合、データには一般に非常に大きな雑音が含まれているものと考えなければならない。つまり、金融商品Aを購入するかどうかというのは、データベースに含まれている項目からは本質的には決定できないはずであるから、もともと非常に精度のよい分類ルールを生成することは期待できない。また、上に挙げた半導体の不良要因分析の場合では、不良の発生がランダムに変動する要因の影響を受けるため、データには大きな雑音が含まれる。このような場合にも、確定的なルールを生成することを望むことが困難な場合が多い。

【0015】上記のような問題の場合には、データの大体の特徴を表現するような分析方法が有効である。しか

しながら、上記従来技術のような方法では、多くの項目の値を組み合わせて少しでも分類精度のよいルールを探しに行くので、一般にIF部にあらわれる条件の数は増大するが、そのルールにあてはまる事例の数は減少するという現象が起こる。これではデータの大体の特徴を理解するという目的を満足することは困難である。

【0016】実際のデータベースには非常に多くの種類の情報が格納されている。上記の半導体品質管理データの例でいえば、ウェハ番号・製造開始年月日など、また顧客データの例でいえば、氏名・電話番号など、明らかに分析の目的とは関係ないものも含まれている。一方で、半導体品質管理データの場合には製品種別コード、顧客データの場合には住所のように、場合によっては分析有効な情報もある。

【0017】そこで、そのように多くの種類の情報からなるデータを用いて上記の例のような分析をする場合には、まずどのような項目のデータを使って分類を試みるかということ、利用者があらかじめ指定しなければならない。項目の種類が増えれば増えるほど、この作業は複雑になる。もしかすると関係があるかもしれないような項目を全て含めて分析しようとすれば、必然的に使用する項目数は増え、処理時間の増大を招く。したがって効率よく分析をしようと思えば、注意深く使用項目を選ばなければならず、利用者のノウハウの程度によって、分析の結果の質が大きく左右されることになる。

【0018】さらに、項目によってはいろいろな見方での分析が必要になる場合がある。例えば、顧客データの例で、地域別に商品購入の動向が異なるような場合、住所という項目を県別で考えるべきか、東北地方/関東地方といったレベルで考えるか、東日本/西日本の2分類で考えるべきかというのは、基本的には分析して見なければ分からない。このような場合を全て尽くそうと思えば、何度も繰り返して従来技術による分析を試みる必要があり、利用者の負担が増大するという問題点がある。

【0019】これを避けるために、いろいろなレベルでの見方を全てデータ項目に追加するという方法もある。住所の例でいえば、県別/地域別/東西別という項目を分析すべきデータ項目として考えることは可能である。しかし、従来技術による分析では、項目間の意味的な関連は一切考慮していないため、処理に大きな無駄が生じる。

【0020】例えば、県別というレベルでの分類を試みている間は、それより上位のレベルである地域別や東西別という項目の値は本来考慮する必要がないが、それにもかかわらず、従来技術では無駄な分析を行ってしまう。さらに、項目の処理順序によっては、「IF 住所が関東地方 AND 住所が神奈川県 THEN …」といった、明らかに意味的に冗長なルールが生成される可能性がある。計算機による分類が目的の場合には、このような冗長性は分類精度には影響しないが、人間がデータの

特徴を理解するという目的には明らかに有害である。

【0021】また、実際のデータベース中には欠損値と呼ばれる未知のデータ項目が含まれていることが多い。統計手法などで分析をするときには、欠損値は単に無視し、データがないものと見なさざるを得ない。上記ルールインダクションの方法でも、欠損値をもつデータ項目は分類精度には影響しないので、分類ルール中には現れることがない。

【0022】しかし、データの値が欠損していること自身に意味がある場合がある。例えば住所という項目が欠損値ということが、匿名で銀行口座を作っていることを意味する場合、このことが金融商品Aの購入に影響を与えることが考えられる。このような場合には、「IF 住所が欠損値 THEN …」というルールが意味を持つことになる。従来の方法では、このようなルールは一般的には生成できず、また生成するためには欠損値を明示的に特定の値に変換するような人手による処理が必要になるという問題点があった。

【0023】さらに、上記従来技術による方法では、あくまで分類精度を優先してルールを生成するため、できるだけ多くの事例を説明するという意味で一般的なルールが必ずしも先に生成されるとは限らない。一方、今考えているようなデータ分析方法では、処理時間が長い場合に途中で利用者が割り込みをかけることもある。このようにときに、利用者にとって利用価値の高い一般的かつ単純なルールから先に生成されていれば、割り込みをかけた時点までに生成されたルールを利用することが可能になるが、従来技術による方法ではこのような利用ができないという問題点もある。

【0024】本発明の目的は、データベースなどに格納されたデータを分析して利用者にとって有用な表現に変換して出力する方法および装置において、上述した問題点を解決することにある。また、本発明の目的は、対象となるデータの量と性質に応じて変換処理方法の一部を自動的に決定することにより、利用者が処理方法に関する特別な知識を有しない場合にも、データの特徴を精度よく表現する結果を生成し、かつ上記の処理を高速に実行し、さらには対話的な計算機環境で上記処理を実行する場合に利用者にとって良好な操作環境を提供するのに適したデータ分析方法および装置を提供することにある。

【0025】

【課題を解決するための手段】本発明は、複数のデータ項目の具体的な値である数値または記号値のデータの集まりを1つのレコードとし、複数のレコードからなるデータを入力して分析することにより、該データの特徴を表現するルールを生成出力する際に、前記複数のデータ項目の中から、ルールの条件節と結論節に用いるデータ項目を選択し、選択したデータ項目が数値の値を持つデータ項目であるときは、数値を記号値に変換し、データ項

目間の相互関係を項目名称と記号値の組みを1つまたは複数の組み合わせたルール形式で表現した候補ルールを複数作成し、該候補ルールが表現する相互関係の強さを評価する評価尺度を算出し、作成した複数の候補ルールのうち該評価尺度が大きな候補ルールを1つまたは複数個求め、求めた候補ルールを出力することを特徴とする。

【0026】また本発明は、複数のデータ項目の具体的な値である数値または記号値のデータの集まりを1つのレコードとし、複数のレコードからなるデータを入力して分析することにより、該データの特徴を表現するルールを生成出力する際に、前記複数のデータ項目の中から、ルールの結論節に用いるデータ項目である結論項目を1つ選択し、前記複数のデータ項目の中から、ルールの条件節に用いるデータ項目である条件項目を選択し、選択したデータ項目が数値の値を持つデータ項目であるときは、数値を記号値に変換し、条件項目の項目名称とその記号値の組みを1つまたは複数の組み合わせで条件節を作成し、結論項目の項目名称とその記号値を組み合わせで結論節を作成し、作成した条件節と結論節を組み合わせで複数の候補ルールを作成し、作成した候補ルールの良さを評価する評価尺度を算出し、作成した複数の候補ルールのうち、算出した評価尺度が大きな候補ルールを1つまたは複数個求め、求めた候補ルールを最終的に生成したルールとして出力することを特徴とする。

【0027】以下、本発明を具体的に説明する。まず、分析の対象となる項目（結論項目と呼ぶ）を指定する手段を設ける。上記従来技術で説明した例の場合は、「金融商品Bの購入履歴」という項目が結論項目となる。それ以外の項目を条件項目と呼ぶ。

【0028】次に、対象とするデータ中に含まれる事例の数を調べ、その数に応じて生成されるルールに含まれる条件節数の最大値、各ルールの最小カバー率、ルール群のカバー率、生成するルール数、を決定する。条件節数とは、ルールのIF部に含まれる述語（「アルミー層配線ショート発生率が5%以上」「年齢が40歳以上」など）の数である。あるルールのカバー率とは、ルールのIF部を満足する事例の数の、全体の事例の数に対する割合である。ルール全体のカバー率とは、生成されたルール群に含まれるルールのいずれかの条件節を満足する事例の数の、全体の事例の数に対する割合である。

【0029】次に、データ中の各項目に関してその値が記号属性を持つか数値属性を持つかを調べる。数値属性を持つ項目に関しては、数値の範囲をあらかじめ定められた方法で複数の区間に分割し、各区間に異なる記号値を割り当てることにより、数値属性を記号属性に変換する手段を設ける。分割する区間の数は、あらかじめ定められた値の範囲で、結論項目との依存関係がもつとも大きくなるような数を選択する手段を設ける。この時に欠損値は、あらかじめ定めた特別の値として表現するもの

とする。

【0030】上記分割する区間には、それぞれメンバシップ関数を割り当てる手段を設け、数値を記号値に変換する時の情報量の損失を小さくする。

【0031】記号属性を持つ項目に関しては、各項目のとり得る記号値の総数を調べる手段を設け、その数があらかじめ決められたしきい値よりも小さい場合にのみ条件項目とし、それ以外の場合にはルール生成には用いないような処理手段を設ける。

【0032】さらに、記号属性を持つ2つの項目の対に関し、記号値の間の依存関係が1対多の関係を持つ場合には、1に対応する項目が多に対応する項目の上位概念であるという性質を記憶しておく手段を設ける。この記号属性間の上位概念/下位概念という性質は、ルール生成に使用するデータをRDBなどのデータベースシステムにおいて複数の関係テーブルから定義する場合には、要素テーブル間の関係定義から自動的に抽出する手段を設けることによって得ることもできる。

【0033】また、数値属性を前述した手段により記号属性に変換した項目に関しても、上記の性質を調べることで、上位概念/下位概念の性質を求めることも可能である。

【0034】次に、条件項目として選択された項目のそれぞれに関して欠損値の有無を調べ、欠損値が存在する場合には欠損値であるかどうか結論項目の値に影響を与えるかどうかを統計的な検定手法などにより調べ、影響があると判断した場合には欠損値をルールの条件部に利用できることを示すあらかじめ定められた別の記号値に変換する手段を設ける。

【0035】上記の処理の結果をディスプレイなどの表示装置に表示し、表示された情報を利用者が必要に応じて修正できる手段を設けておく。

【0036】ルール生成としては、カバー率とヒット率から決まるルールの評価尺度を定義し、最大条件節数以下の条件節をもち、かつカバー率が所定の値以上の値をもつ全てのルール候補の中で、大きな評価尺度をもつルールから上記の手段により定められたルール数を選択し、結果として出力装置に出力する手段を設ける。ここで、ヒット率とは、(ルールの条件節と結論節の両方を満たすレコードの数) ÷ (ルールの条件節を満たすレコードの数) であるものとする。または、評価尺度の大きな順に候補ルールを選択していき、ルール群としてのカバー率が所定の値以上になるまでルールを選択し、生成ルール群とする手段を設ける。

【0037】前記評価尺度は、ルールのカバー率とヒット率と条件節数から算出する手段を設けておくこともできる。

【0038】数値属性を記号属性に変換するときにメンバシップ関数を用いた場合には、カバー率およびヒット率は各事例のメンバシップ値を用いて算出する手段を設

けておく。

【0039】上記ルール生成では、条件節数の小さな候補ルールから順に生成し、生成済みの候補ルールに新しい条件節を追加することにより新しい候補ルールを生成する手段を設ける。さらに、生成した候補ルール数に応じてルール生成処理の進行状況を表示装置に表示する手段を設ける。また、利用者が処理の中断を指定する手段を設けるとともに、処理の中断が指定された場合には、その時点までで生成された候補ルールの全部または一部を出力装置に出力する手段を設けておく。

【0040】ある候補ルールに対して条件節を追加した場合にとりうるルール尺度の最大値を算出する手段を設け、その最大値に応じて条件節を追加するかどうかを判定する手段を設ける。

【0041】ルール生成手段の中間結果または最終結果として得られたルール群に関し、ルール群全体、または1つまたは複数の選択されたルールに関するルール尺度を最大化するように、数値属性を記号値に変換する際に設けたメンバシップ関数の形状を変更する手段を設ける。

【0042】

【作用】本発明によれば、対象とする事例の集合から、条件節数の最大値・各ルールの最小カバー率・ルール群のカバー率・ルールの条件部に使用される条件項目・数値属性から記号属性への変換方法・記号属性間の上位概念/下位概念関係・欠損値の扱い方法など、データ分析処理のパラメータを全て自動的に決定することができる。

【0043】したがって、データを分析しようとする利用者は、説明しようとする結論項目さえ指定すれば、あとは計算機処理によってルール形式で表現されたデータの特徴を結果として得ることができるため、分析作業の手間を軽減し、また利用者のノウハウの存在を前提とすることなく、精度のよいデータ分析方法を提供することが可能になる。もちろん、自動的に決定された上記パラメータを表示装置に表示し、利用者が確認し必要に応じて修正する手段を設けることにより、利用者のデータ分析の目的に、より適した処理を行うことが可能である。

【0044】また、ルールの良さを評価する尺度として、カバー率・ヒット率・条件節数から決まる評価尺度を用いることにより、利用者がデータの特徴を理解するのに適したルールが生成される。

【0045】数値属性を記号属性に変換する際にメンバシップ関数を用いる方法では、例えば同一の区間に属する値でも、その区間内の中心近くの値を持つのか、区間の端に近い値を持つのかメンバシップ関数の値として区別できるため、記号化することによる情報量の損失を抑えることができ、生成されるルールの精度が向上し、かつ生成されるルールとしては記号値で表現された理解しやすいルールを生成することができる。



【0046】また、候補となるルールの生成を条件節数の小さなルールから順に行うため、事例の数や項目数が非常に大きくて処理が利用者の想定した時間内に終了せず、利用者の指示によって処理が中断された場合でも、その時点までに生成されたルールがデータの特徴の概要を理解するのに有効なものになる。

【0047】また、順次候補ルールを生成する際に、条件節を追加して生成されるルールの尺度の最大値をあらかじめ評価しておくことにより、無駄な候補ルールの生成を避けることができ、効率的な処理が実現できる。結果として、処理を高速化することができるようになる。

【0048】さらに、途中結果または最終結果として得られた（候補）ルール群に含まれるメンバシップ関数で表現された記号値を、ルール尺度が最大になるようにメンバシップ関数の形状を調節することによって、ルールに含まれる記号の意味を、データの特徴をもっともよく説明するものとすることができる。例えば、当初の処理で年齢という項目に対して「30から39」という区間を設定したとし、金融商品の購入する顧客が何らかの理由で33歳から37歳の区間に集中した場合でも、その特定の区間を自動的にメンバシップ関数で表現できるようになる。つまり、その区間が特徴のある顧客セグメントの存在を示すという情報を利用者に対して提示することができるようになる。

【0049】

【実施例】以下、図面を用いて本発明の実施例を説明する。本実施例では、前述の金融商品購入顧客の分析の例を用いる。

【0050】図2は、分析の対象とするデータの例を示す。ある金融商品の購入動向に関して動向の分析をすることを考える。図2のデータには、顧客番号・氏名・年齢・性別・住所・支店コード・県コード・地域コード・預金残高・ローン残高・クレジットカード種別などの情報とともに、ある金融商品の購入実績が格納されている。

【0051】このデータを分析しようとする利用者は、顧客のどんな属性が金融商品の購入に影響しているかを分析することにより、ダイレクトメールや訪問販売など、その金融商品拡販のための戦略（どのような顧客にアプローチをするのが最適か、など）を検討することを目的としている。

【0052】図3に、分析処理の全体フローを示す。まず処理301では、分析対象データを定義し、図2のような分析対象データの表を作成する。次に処理302で、分析対象データの読み込みを行う。

【0053】処理303では、まず「商品購入」という項目が結論項目であることを指定する。次に、データの各項目を調べ、項目が記号属性を持つか数値属性を持つかを調べる。各項目について、条件項目として選択するかどうかを判定し、選択された項目が数値属性をもつ項

目である場合には、それを記号値に変換する。最終的に、ルール生成処理の入力となるデータを生成する。

【0054】処理304では、事例数などの情報に基づいて、ルール生成処理のパラメータ（最大条件節数、最小カバー率、など）を決定する。処理305では、決定したパラメータを表示装置に表示し、必要ならば利用者がパラメータを修正する。処理306ではルール生成処理を実行し、処理307で結果を出力装置に出力する。

【0055】なお、図3ではルールの学習処理は省略したが、処理307の後、利用者の指示により随時学習処理を実行することができる。

【0056】図4に、以上の処理を実行するためのシステムのブロック構成図を示す。システムは大きく3つの処理モジュールからなる。利用者とのインタラクションを制御するユーザインタフェースモジュール401、分析対象のデータからルールを生成するルール生成処理モジュール402、およびメンバシップ関数の学習によって生成されたルールの微調整を実行するルール学習処理実行モジュール403、である。

【0057】モジュール401は、利用者からの指示をコマンドとして認識し、指示された方法に基づいて分析対象データベース404からデータを入力し、ルール生成処理実行モジュール402の入力となるルール生成用ワークファイル405を生成し、同時にルール生成処理を制御するためのパラメータ類の情報をルール生成パラメータファイル406に出力する。利用者からのルール生成実行コマンドがモジュール401に入力されると、モジュール401はモジュール402を起動する。

【0058】モジュール402が起動されると、ファイル405とファイル406の情報を読み込み、ファイル405のデータに対して、ファイル406で指定されたパラメータでルール生成処理を実行し、結果をルールファイル407に出力する。ルールファイル407には、生成されたルールの情報と、生成されたルールに含まれる述語を定義するメンバシップ関数の情報が格納されている。

【0059】モジュール402の実行が終了すると、制御はモジュール401に戻り、生成されたルールをあらかじめ定められた形式で利用者に表示する。利用者は、表示結果に基づいて、以上の処理をパラメータを変更して再度実行するなどの作業を続行することができる。

【0060】ここで生成されたルールの学習コマンドがモジュール401に入力された場合、モジュール401は、モジュール403を起動する。モジュール403が起動されると、まずデータベース404とルールファイル407の内容を読み込み、ルールファイル407に含まれるそれぞれのルールに対してメンバシップ関数の学習を実行する。処理終了後、学習したメンバシップ関数の情報を再びルールファイル407に格納した後、制御をモジュール401に戻す。

【0061】図3に示した処理のうち、処理306がモジュール402で実行され、残りの処理はモジュール401で実行される。

【0062】以下、各処理の詳細を、図3のフローにしたがって順に説明する。

【0063】図5および図6は、処理301において実行される分析対象データ定義処理の内容を説明する図である。分析の対象となるデータは、複数のリレーションテーブルに格納されており、それぞれ図5(a)の顧客情報テーブル、図5(b)の支店情報テーブル、図5(c)の商品購入履歴テーブルからなる。図6は、図5の3つのテーブルの相互の関連を指定し、分析対象となる図2のテーブル構造定義を模式的に示した図である。これらの定義は、リレーショナルデータベース管理システム(RDBMS)で提供される機能を用いて行うことができる。

【0064】図3の処理302では、図2で示した分析対象データをデータベースから読み込み、計算機内のメモリに展開する。データ読み込み時には、各項目の値が数値か記号値かということを判定し、かつ数値の場合には最大値/最小値、記号値の場合には異なる記号の数を各項目毎に求めておく。

【0065】図7に、図3の処理303における条件/結論項目の指定画面を示す。処理303では、この図7のように画面表示を行い、ユーザに条件項目と結論項目の指定をさせる。

【0066】まず、データ読み込み直後にはいずれの項目も選択されておらず、図2の分析対象データの各項目の名称が未使用項目名リストボックス701に表示されており、条件項目名リストボックス702、および結論項目名テキストボックス703の内容は空である。

【0067】ここで、リストボックス701の表示の際には、処理302で求めた記号値をもつ属性のうち、異なる記号の数があらかじめ定められた数、例えば20個以上の場合には条件項目/結論項目として選択することを不可とし、その項目名が選択不可であることを表示するものとする。これは、記号値の数が事例数に比べて大きな項目を条件または結論項目として指定すると、生成されるルールが細分化されて利用価値が低下することを防ぐためである。例えば、図7の顧客番号・氏名・住所が相当する。図7では、項目名の先頭に「\*」を付加してその項目が選択不可であることを表現しているが、もちろん網掛け表示/淡色表示などを用いても同様の効果が実現できる。このようにすることにより、ルール生成に使用しても意味のない項目を使用することを防止することができる。

【0068】リストボックス701中の(選択可能な)項目をマウス等の入力装置を用いて複数選び、ボタン706をクリックすることにより、選択中の項目が条件項目として選択され、リストボックス701からその項目名が削除され、リストボックス702にその項目名が追

加されて表示される。同様に、リストボックス701の項目を1つ選択してボタン709をクリックすることにより、結論項目が選択される。リストボックス702および703の項目を削除するときも同様に、削除項目を選択してボタン707、710をクリックする。

【0069】ボタン705をクリックすることにより、リストボックス701に表示されている項目で選択不可の項目以外を全て条件項目に追加することができる。ボタン708は、逆にリストボックス702中の全ての項目を条件項目から削除するためのボタンである。

【0070】これらの選択が終了した後、ボタン711をクリックすると、条件/結論項目選択処理が終了する。また、ボタン712をクリックすると、その時点での項目選択情報を全て破棄する。

【0071】次に、ボタン704をクリックしたときに実行される条件項目自動選択処理について説明する。ただし、この処理は結論項目が指定されているときに実行できるため、ボタン704は結論項目が指定されているときだけ有効であるものとする。

【0072】図8は、条件項目自動選択処理の処理フローである。図7において、結論項目(テキストボックス703)の指定の後、ボタン704をクリックすると図8の処理が実行開始する。

【0073】処理801では、使用不可の項目を除く全ての項目に関して処理802、803を繰り返す。処理802では、現在処理している項目が数値属性を持つか記号値属性を持つかを調べ、数値属性を持つ場合には処理803を実行する。

【0074】処理803では、処理中の項目の値をあらかじめ定められた個数、例えば5個の記号値に変換する。図2中の「年齢」という項目の例では、項目値のとりうる値の範囲を例えば「20未満」「20以上30未満」「30以上40未満」「40以上50未満」「50以上」の5つの範囲に分割し、5つのカテゴリに分類することになる。この分類の方法には、いろいろな方法があるが、例えば、下記の3つの方法などがある。

【0075】1)等範囲分割:とりうる値の範囲を等分割する。例えば年齢の最小値が0、最大値が75とすれば、0~14、15~29、30~44、45~59、60~75と分類する。

【0076】2)平均/標準偏差分割:年齢の値の分布の平均と標準偏差を求め、その値を元に分割する。例えば平均を $\mu$ 、標準偏差を $\sigma$ とすると、 $\mu - 0.84\sigma$ 、 $\mu - 0.25\sigma$ 、 $\mu + 0.25\sigma$ 、 $\mu + 0.84\sigma$ という値を分割点とする。この値は、分布が正規分布をすると仮定したときに、各分類に含まれる確率を等しくする値である。

【0077】3)等数分割:実際の年齢の値の分布を調べ、各分類に含まれる事例の数が等しくなるように分類する。

【0078】どの方法を用いて分類するかは、あらかじめ

めデフォルトの方法を用意しておき、必要に応じてユーザが指定できる手段を用意しておくものとする。

【0079】上記の方法によって分類された数値データは、便宜的に「カテゴリ1」「カテゴリ2」～「カテゴリ5」という値を持つ記号値データに変換される。

【0080】処理804では、結論項目として指定された項目以外の項目中の任意の2つの項目の組からなるすべての組み合わせについて、2つの項目間の類似度尺度をすべて計算する。類似度尺度は、次のように計算する。

【0081】まず、2つの項目をそれぞれX、Yとし、\*

$$I(X; Y) = \sum_i \sum_j P(x_i, y_j) \log [P(x_i, y_j) / P(x_i) P(y_j)] \quad (1)$$

【0083】ここで、 $\sum_i$ は添字iに関する和、 $P(x_i)$ は $X=x_i$ となる確率、 $P(x_i, y_j)$ は $X=x_i$ 、 $Y=y_j$ となる確率を表す。確率はデータ中のレコードの出現頻度から計算するものとし、例えば、 $P(x_i) = N(x_i) / N$ 、 $N$ は全レコード数、 $N(x_i)$ は $X=x_i$ となるレコード数とする。

【0084】 $I(X; Y) \geq 0$ という性質があり、 $I(X; Y)$ の値が大きければ大きいほどXとYの間には関連があることを示しているとみなせる。極端な場合、XとYが全く独立の場合には、 $P(x_i, y_j) = P(x_i) P(y_j)$ となり、 $I(X; Y) = 0$ となる。

【0085】この性質から、項目Xと項目Yの類似度尺度 $D(X; Y)$ を、次の式で定義することができる。

$$D(X; Y) = 1 / I(X; Y) \quad (2)$$

ただし、 $I(X; Y) = 0$ のときは、Dの値は非常に大きな固定値をとるものと約束しておく。

【0086】処理805では、処理804で結論項目以外のすべての項目の組み合わせごとに計算した類似度尺度Dを基準にして、結論項目以外の項目をクラスタリングする。すなわち、 $D(X; Y)$ を項目Xと項目Yの距離とみなし、全ての項目間の距離に基づいて、距離の近い項目をまとめあげ、いくつかのクラスタとして分類する。

【0087】クラスタリングの方法には様々なアルゴリズムが知られており、その詳細は一般的な統計学の教科書・文献に述べられているので、ここでは説明を省略する。例えば、k-means法と呼ばれる方法を用いるものとする。これによって、条件項目の候補となる項目の内、お互いに類似度（関連度）の高い項目を一つのクラスタにまとめあげることができる。

【0088】多くのアルゴリズムではクラスタ数をあらかじめ指定するが、ここではデフォルト値として全データの項目数の $1/10$ をクラスタ数とするものとする。すなわち、図2の分析対象データの項目数が62項目の場合には、 $62 \div 10 = 6$ （切り捨て）個のクラスタを生成するものとする。もちろん必要に応じてユーザがクラスタ数を直接指定する手段も用意しておくものとする。

【0089】処理806では、結論項目とそれ以外の項

\*そのとりうる値を $x_i (i=1, N_x), y_j (j=1, N_y)$ とする。X、Yは、もともと記号値を持つ項目であっても、上記の方法により数値を記号値に変換した項目であってもよい。XとYを確率変数と考えると、その間に相互情報量が定義される。

【0082】相互情報量に関しては、多くの情報理論の教科書に解説されている。相互情報量とは、例えば $X=x_i$ という情報が与えられたときに、Yの値の発生する確率分布が持つ情報量がどう変化するかに関係した量であり、いわばXとYの間の関連度を示す量である。相互情報量 $I(X; Y)$ は、次の式で定義される。

目の組み合わせに関して、依存度尺度を計算する。依存度尺度は、例えば処理804で説明した相互情報量 $I(X; Y)$ を用いることができる。

【0090】処理807では、処理805で分類されたクラスタについて、それぞれのクラスタに属する項目の中で最も結論項目との依存度尺度が大きな項目を選択し、それらの項目を条件項目として採用する。

【0091】以上、図8の処理により、数値・記号値をもつ項目が混在したデータから、結論項目との関連度が最も高く（依存度が大きい）、かつ条件項目間の関連度をできるだけ低く（独立性を大きく）した条件項目を自動的に選択することが可能になる。これによって、以下の処理で生成されるルールの高精度化が図れ、かつ類似の条件項目を使用しないので、ルール間の独立性を高めることができるので、生成ルールの可読性／有用性を高めることができるという効果がある。

【0092】処理303においては、上述の処理で条件項目に選択された項目に対して、欠損値対策処理を実行する。欠損値対策処理は、次のステップを条件項目に選択された各項目に対して行う。

【0093】ステップ1：選択中の項目に対し、欠損値データが含まれるかどうかを調べる。欠損値データは、特別な数値（例えば負の最大値）または記号値（例えば空の文字列）で表現されているものとする。欠損値が含まれていなければ、処理を終了する。

【0094】ステップ2：選択中の項目の値が欠損値か欠損値以外かという分類をしたときの、結論項目Zとの相互情報量 $I(X; Z)$ を求める。Xは、選択中の項目の値が欠損値か欠損値以外かという値を持つ確率変数である。

【0095】ステップ3：ステップ2で求めた相互情報量 $I(X; Z)$ があらかじめ定められたしきい値 $I_{th}$ よりも大きな場合には、欠損値であるかないかということと結論項目Zとが所定の程度以上に関連しているということだから、選択中の項目の欠損値を新しい記号値として追加する。

【0096】以上の処理により、欠損値であること自身に意味があるとステップ3で判断された場合には、「欠

損値」という特別の記号値がその項目に追加され、後述のルール生成にその記号値を使用することができるようになる。それ以外の場合には、その項目の欠損値は、後述するルール生成処理では一切無視される（どの記号値にも属さず、生成ルールにも決して現れることがない）特別な値として扱われる。

【0097】図3の処理303における数値→記号値変換処理は以下のように行う。

【0098】図9は、数値→記号値変換方法の一覧表示である。ユーザは、表示装置場の図9のような表示を参照して、マウスなどを用いて各項目について変換方法を指定する。

【0099】図9において、項目名表示欄901には、本来数値属性をもち、記号値に変換される項目の名称が表示される。記号数表示欄902には、記号値に変換した後の異なる記号の数が表示される。記号名称表示欄903には、記号値の名称が表示される。

【0100】図9の例では、年齢は20未満、20代、…、預金残高は少ない、ふつう、多い、ローン残高は小さい、中くらい、…、とそれぞれ異なる記号値が割り当てられている。変換方法表示欄904には、図8の処理803で説明した数値→記号値変換方法が表示される。図9の例で、年齢に対応する変換方法がユーザ指定となっているが、このユーザ指定とは、ユーザが分割方法を直接指定したことを示している。

【0101】各項目に対する記号値変換方法の指定が終了した後、ボタン906をクリックすることで変換方法指定処理が終了し、ボタン907をクリックすることでそれまでの指定内容が破棄される。ボタン905をクリックすると、図8の処理803でも説明したようにデフォルトの記号値変換方法で全ての項目の変換方法をリセットする。

【0102】各項目毎に変換方法を指定するには、欄901の項目名の部分をダブルクリックする。これにより、図10の変換方法指定画面が表示され、当該項目の変換方法を指定することができるようになる。

【0103】図10において、オプションボタン1001は、選択中の項目の記号値変換方法として、等数分割・等範囲分割・平均／標準偏差分割・ユーザ指定分割のどの方法を用いるか、およびそれぞれの分割方法をクリ

スプ的に行うかファジィ的に行うかを指定する。

【0104】テキストボックス1002には、記号値に変換する際の記号数を入力する。

【0105】記号名称表示欄1003には、現在の記号名称が表示されている。ユーザが、この記号名称を変更したい場合には、直接、欄1003の該当領域をクリックすることにより記号名称を入力することができる。これは、各記号値の最小値表示欄1004および最大値表示欄1005に関しても同様である。ただし、最小値または最大値を直接ユーザが変更した場合には、オプショ

ンボタン1001は自動的にユーザ指定分割方法を指定した状態に変化するものとする。

【0106】欄1004または欄1005の値が変更されると、レコード数表示欄1006の内容もそれに応じて自動的に変更される。これは、メモリ中に記憶されているデータを用い、それぞれの記号値に対応するレコードが何件あるかをカウントすることによって行う。

【0107】オプションボタン1001のうちクリスプ分割またはファジィ分割のいずれかを指定することによって、欄1004および欄1005で表示される最小値／最大値の範囲で記号値変換処理をクリスプ的に行うかファジィ的に行うかを指定できる。ファジィ的に行うことにより、記号値変換によって失われる情報量を最小限にし、生成されるルールの精度を向上することができる。以下、この処理について説明する。

【0108】ファジィ的な分割をする際にも、各記号値の値の範囲は図10の欄1004および欄1005には最小値／最大値として表示されるが、それをファジィとして解釈する。図11は、このようなファジィ的な分割を図10のローン残高の記号値変換の例で説明する図である。

【0109】図11の横軸はローン残高の数値を万円単位で表示し、縦軸は「小さい」「中くらい」「大きい」という（ファジィ）記号値の適合度を表す。「小さい」という記号値にはメンバシップ関数1101が、「中くらい」にはメンバシップ関数1102、「大きい」にはメンバシップ関数1103が対応する。メンバシップ関数の形状は、傾きの絶対値が一定、各記号値の分割点（図では、「小さい」「中くらい」の分割点である500、「中くらい」「大きい」の分割点である1000の2カ所）でそれぞれのメンバシップ関数値が0.5になり、かつどのような値に対しても0以外の値をもつメンバシップ関数は高々2個であり、かつメンバシップ関数値の和が1.0になるように設定することができる。

【0110】例えば図11で、「ローン残高＝600（万円）」という数値は、「中くらい」という記号値に対してが0.8、「大きい」という記号値に対して0.2というメンバシップ値を持つことになる。すなわち、「ローン残高＝600」という値をもつレコードは、「ローン残高」が「中くらい」というレコードが0.8個分、「大きい」というレコードが0.2個分の2つに分割されたように振る舞うものとみなすことができる。

【0111】このようなファジィ分割の場合には、図10の欄1006に表示されるレコード数も、それぞれのメンバシップ値の合計と考えることができ、一般に実数の値をとることになる。

【0112】次に、図3の処理304におけるルール生成パラメータ算出処理を説明する。ルール生成処理を制御するパラメータとしては、次の5種類のパラメータを

- ・生成ルール数
- ・全体カバー率
- ・最大条件節数
- ・最小カバー率
- ・カバー率優先係数

【0113】生成ルール数は、ルール生成処理で生成すべきルール数を指定するパラメータである。デフォルト値は、使用する条件項目数Nから次の関係式にしたがって定める。

【0114】(生成ルール数のデフォルト値) =  $N \times$  (最大条件節数)

これは、一般的に条件項目数が多くなればデータの特徴を説明するために生成すべきルール数も増えることを反映させたものである。

【0115】カバー率とは、あるルールの条件部を満たす事例の数の、全体の事例数に対する割合と定義する。全体カバー率は、生成したルールに対して、少なくとも1つのルールの条件部を満たす事例の数の、全体の事例数に対する割合を指定するパラメータである。生成するルール数を大きくすれば全体カバー率は大きくなるという関係がある。全体カバー率に対してはデフォルト値は用意しない(すなわちユーザの明示的な指定によってのみ値を設定できる)。

【0116】最大条件節数は、ルールの条件部に含まれる述語数の最大値を指定するパラメータである。例えば、最大条件節数3に指定した場合、

もし、X1がA1 かつ X2がA2 かつ X3がA3 ならば ZはB

というルールは生成されるが、さらに「～ かつ X4がA4 ならば ～」というルールは(条件節数が4となるので)生成されないことになる。デフォルト値は3とする。

【0117】最小カバー率は、生成されるルールのカバー率の最小値を指定するパラメータで、ここで指定された値以下のカバー率をもつルールは、特殊なルールとして生成されない。最小カバー率のデフォルト値は、 $100 / (\text{生成ルール数})$  (単位パーセント)とする。これは、一般的に生成ルール数が多いほど細かなルールでも必要とされることを考慮したものである。

【0118】カバー率優先係数は、ルールに対する評価尺度におけるカバー率と精度の関係を指定するもので、詳細は後述のルール生成処理において説明する。デフォルト値は1.0(カバー率を最大限に優先)とする。

【0119】図3の処理305では、処理304で算出\*

$$\mu(A \rightarrow B) = P(A) \cdot \beta \cdot \log[P(B|A) / P(B)] \quad (3)$$

【0129】ここで  $a^b$  はaのb乗を意味する。P

(A)は分析対象とするデータの中で条件部Aが満足される確率、すなわち事例全体の中でAという条件を満たす事例の割合、を表す。同様に、P(B)は結論部Bが満足される確率、P(B|A)はAという条件を満たす

\*したルール生成パラメータを表示装置に表示し、必要に応じてユーザが修正できるようにする。図12に、表示画面例を示す。

【0120】テキストボックス1201～1205には、図3の処理304で算出されたルール生成パラメータの値が表示される。図の例では、生成ルール数=24、全体カバー率=指定なし、最大条件節数=3、最小カバー率=4%、カバー率優先係数=1.0である。

【0121】ユーザは、テキストボックスを直接クリックすることで表示されたパラメータの値を入力することができる。図12の表示画面には、分析対象となるデータの他の情報(データファイル名、レコード数、条件項目数、結論項目とその記号値)が表示領域1206に表示されているので、ユーザはこれらの情報を参考にしながらルール生成パラメータの値を修正することが可能である。

【0122】パラメータの修正後、ボタン1207をクリックすることにより修正した値が有効になり、ボタン1208がクリックされると修正した値は破棄される。

【0123】図3の処理306では、処理304および処理305で指定されたルール生成パラメータにしたがってルール生成処理を実行する。

【0124】ルール生成処理では次の形式のルールを生成する。

・もし X1 が A1 かつ X2 が A2 かつ ...  
かつ Xn が Anならば Y が B

ここで、Xiはそれぞれ条件項目名称、Aiは項目Xiの記号値の名称である。同様に、Yは結論項目名称、Bは項目Yの記号値の名称である。

【0125】「XiがAi」という組を条件節と呼ぶ。上記のルールは、条件節をn個もつルールの例である。条件節全体を条件部とよぶ。また「YがB」を結論節と呼ぶ。本実施例では結論部は1つの結論節からなるものとしておく。

【0126】ルールには、評価尺度と呼ばれる実数値を割り当てる。評価尺度の大きなルールほど、ルールとしての価値が高いものとみなす。

【0127】図3の処理306では、条件節と結論節の組み合わせの中から、ルール評価尺度の大きなルールを選び出すという処理を実行する。これにより、価値が高いルールを選ぶことができる。

【0128】本実施例において、「もし A ならば B」というルールの評価尺度 $\mu(A \rightarrow B)$ は次式のように定義される。

という前提で結論部Bが満足される確率を表す。 $\beta$ は図3の処理304で説明したカバー率優先係数である。式(3)を事例の数で書き直すと、 $P(B|A) = P(A \& B) / P(A)$ であることから、

$$\mu(A \rightarrow B) = [N(A)/N]^\beta * \log[N \cdot N(A \&B)/N(A)N(B)] \quad (4)$$

となる。

【0130】ただし、

N: 全体の事例数

N(A): 条件部Aを満たす事例の数

N(B): 結論部Bを満たす事例の数

N(A&B): 条件部Aと結論部Bを同時に満たす事例の数である。

【0131】式(4)の定義の第一因子はルールのカバー率の $\beta$ 乗であり、第二因子は精度に対応する。すなわち、カバー率が大きく(できるだけ広い範囲の事例を説明し)、かつ精度がよい(ルール全体が成立する確率が大きい)ほど、高い評価尺度をもつことになる。一般には広い範囲をカバーするほど精度は低下する傾向にあるので、これら2つの因子は互いに背反する傾向にある。ルール生成パラメータの一つであるカバー率優先係数 $\beta$ は、評価尺度のなかでどの程度カバー率を重視するかを決めるもので、 $\beta=1.0$ の時にはカバー率を最も優先して考慮し、 $\beta=0.0$ の時にはカバー率は考慮せず精度だけを考慮してルールを評価することになる。

【0132】もちろんここで定義したルール評価尺度は、あくまでも1つの候補であり、いくつかの評価尺度の定義を追加すること、またそれらの定義をユーザオプションとして切り替えて使用すること、などが可能である。

【0133】ルール生成処理では、式(4)で定義したルール評価尺度の大きなルールから指定された個数のルールを生成する。探索すべき空間は、すべての可能な条件節の組み合わせになる。

【0134】条件項目10個、それぞれの記号数3程度の問題でも、この探索空間は膨大な広さになる(およそ30の10乗のオーダー)ため、全数探索を現実的な時間で実行することは不可能である。したがって、以下に説明する方法を用いて、探索範囲の制限と枝刈りにより探索を効率化する。

【0135】まず、最大条件節数と最小カバー率を用いてルール探索範囲が制限できる。

【0136】これらの制限後の空間での探索方法として、本実施例では深さ優先探索法を用いる。

【0137】図13に、探索空間のイメージを示す。基本的には条件節述語の組み合わせであるから、木構造の探索手法が使える。

【0138】図13において、例えば、ノード1301に対応するルールは、if X1 is A then ...であり、ノード1302に対応するルールは、if X1 is A and X2 is B then ...である。すなわち、子ノードに対応するルールの条件節は、親ノードに対応するルールの条件節に1つ述語を追加したものになる。各ノードの近傍に付された数値は、当該ノードに対応するルールの条件節の述語を示すものとする。例えば、ノード130

1の近傍には{1}と記載されているが、これは対応するルール中の述語「X1 is A」を示している。また、ノード1302の近傍には{12}と記載されているが、これは対応するルール中の述語「X1 is A and X2 is B」を示している。

【0139】図13では、最大条件節述語数3の場合を示した。木の高さ(ルートから終端ノードまでの最大リンク数)は最大条件節述語数になる。

10 【0140】条件節に現れる述語の順序には意味がないので、例えばノード1303の子ノードは{21}でなく、{23}から始まっていることに注意しておく。また1つの条件節に、同じ条件項目で異なる記号値をもつ2つの述語は含まれない。親ノードの条件節に1つ述語を追加して子ノードを生成する時には、この条件を考慮しておく必要がある。

【0141】深さ優先探索では、図13に示すすべてのノードを次の順序で生成し、評価する。

1→12→123→124→...→12n→13→134→135→...→(n-2)(n-1)n→(n-1)→(n-1)n→n

20 【0142】深さ優先探索の場合は、常にそれまでで評価尺度の最も大きなルールN個(Nは生成ルール数)を保持しながら上記の順序で候補ルールを調べていけばよい。

【0143】次に、上述したルールの探索を高速化するための枝刈り方法について説明する。図13のノード数は前述の例で30の3乗程度であるから、直接全数探索することも不可能ではないが、数10~数100万件のデータを対象にすること、項目数も100程度になる場合も少なくないことから、できるだけ効率のよい探索方法を採用する必要がある。

【0144】効率向上の基本的な考え方は「枝刈り」である。すなわち、ある中間ノードでの状態を評価することによって、その子ノードの評価をせずにすませることである。本実施例におけるルール生成処理の場合、カバー率による枝刈りとルール評価尺度による枝刈りを行う。

【0145】カバー率による枝刈り方法の基本的な考え方は、あるノードに対応するルールを生成した時、そのルールのカバー率がユーザ指定のカバー率下限を下回った場合には、そのノードおよびすべての子ノードを探索する必要がない、ということである。子ノードに対応するルールの条件節は親ノードに対応するルールの条件節に1つ述語を追加したものになるから、親ノードのルールのカバー率が下限を下回ったときは、子ノードのルールのカバー率も必ずその下限を下回ることになる。

【0146】ルール評価尺度による枝刈りについて説明する。探索においてN個のルールを生成する場合、探索のある時点でのN番目の候補ルールの評価尺度を $\mu_{th}$ とする。枝刈りの基本的な考え方は、あるノードをルート

とする部分木におけるルール評価尺度の最大値を求め、その値が $\mu$ thよりも小さな場合には部分木のノードは生成しない、ということである。ここで、いかに部分木に含まれるノードの評価尺度の最大値を推定するかが、枝刈りの効率を決める重要な因子になる。

【0147】以下、この評価方法について説明する。簡\*

$$\mu(A \rightarrow B) = [N(A)/N] * \log[N \cdot N(A \& B)/N(A)N(B)] \quad (5)$$

【0149】条件節Aにいくつかの述語を追加した条件節をA'とすると、 $\mu(A' \rightarrow B)$ の最大値は $N(A')=N(A \& B_i)$ のときに限られる。なぜなら、上式でN、N(B)は定数とみなしてよく、また上式はN(A&B)に関して単調増加であるから、N(A&B)も定数とみなしてよい(すなわ※

$$\mu(A' \rightarrow B) = [N(A')/N] * \log[N \cdot N(A' \& B)/N(A')N(B)] \quad (5)$$

の最大値を求めればよい。これは、上式をN(A')で微分すると必ず0になることから、 $N(A')=N(A \& B_i)$ のときに最大値 $\mu_{\max} = [N(A \& B)/N] * \log[N/N(B)]$ となる。この値が $\mu$ thよりも小さければ、そのノードの子孫ノードを探索する必要がない。

【0150】実際の探索においては、早い段階で評価尺度の高いルールを見つけておくことが、効率向上の点で有効である。早く見つければ、それだけ $\mu$ thの値が大きくなり、枝刈りが実行される可能性を増やすことができる。

【0151】そのための探索戦略の一方法として、図13の木構造をまともに深さ優先で探索するのではなく、まず条件節述語数1の場合をすべて調べ、条件節述語数 $\geq 2$ の場合を深さ優先探索にすることが可能である。すなわち、 $\beta > 0$ ならばカバー率の大きなルールほど大きな評価尺度をもつから、条件節述語数1の場合を先に調べることで、 $\mu$ thの値を早く大きくすることができる。

【0152】このような探索を行うことによって、ルール生成処理を高速化できるだけでなく、ルール生成処理の途中でユーザが処理の中断を指示したときに、中間結果として得られるルールがデータの特徴を全体として表現したルールになる、という効果もある。

【0153】また、生成されたルールはユーザにとってできるだけ理解容易なものであることが望ましい。言い換えると、ルール評価尺度に差がなければ、条件節の述語数は少ないほうがよい。これをルール生成に反映させる方法としては、「条件節述語を追加したルールの精度は、追加する前のルールの精度よりも大きくなければならない」という制約を課すことにする。すなわち、条件節述語を追加してルールを複雑にしたのだから、その分ルールの精度は上がらなくてはならない、ということである。この制約をつけ加えてルール生成を行うことにより、よりわかりやすく利用価値の高いルールを生成することができる。

【0154】さらに、本実施例のルール生成処理では、もともと記号値属性をもつ条件項目間の1対多の依存関係を考慮してルールを生成する。例えば、図2の県コー

\*単のため、カバー率優先係数 $\beta=1.0$ の場合で説明する。一般の値の場合での拡張は容易にできるので、ここでは省略する。

【0148】あるノードにおけるルールを「もしAならばB」とすると、その評価尺度は次式で定義される。

※ち、精度をあえて下げるようなルールは決して評価尺度を最大にしない。したがって、式(5)は下記のように書くことができ、この式について、 $N(A \& B) \leq N(A') \leq N(A)$ の範囲で、

ドと支店コードを考える。支店コードは各支店に固有のコードであるから、支店コードを決めれば県コードが唯一に決まる。一つの県に複数の営業店が存在する場合には、逆は成立しない。このとき、県コードと支店コードの間には1対多の関係が存在するという。

【0155】この2つの条件項目を含むデータに対してルール生成を行うと、支店コードと県コードを同時に条件節述語に含むルールが生成される可能性がある。少なくとも、木構造の探索では候補ルールとして評価尺度が計算される。しかし、これは全くの無駄である。支店コードが条件節述語に含まれていれば、そのコードから県コードは唯一に決まるのだから、県コードに関する条件節述語を追加する必要はない。

【0156】カテゴリカルデータ間の1対多の関係は、分類概念の上下関係と見なすこともできる。すなわち、県コードが上位の分類概念、支店コードが下位の分類概念になる。

【0157】この上下関係は木構造では表現できない。支店コードの例で、さらに地方コードと支店規模コード(大規模・中規模・小規模など)を考える。地方コードは県コードのさらに上位の概念であるから、営業店コードと地方コードの間にも上下関係が存在する。さらに、営業店コードを指定すれば支店規模コードも決まるから、ここにも上下関係が存在する。すなわち、県コードと支店規模コードはどちらも支店コードの上位概念である。逆に、1つのカテゴリカルデータが複数の下位概念を含むこともある(例えば、県コードに対する支店コードと取引先コード、など)。

【0158】本実施例では各条件項目に対して、その上位概念に相当する項目のリストを作成しておく。このリストは、例えば図6に示した分析対象データの定義から自動的に作成することができる。すなわち、リレーショナルデータベースにおいては各テーブルの主キーはテーブル内のレコードを一意的に指定できるものでなければならないから、テーブル内の主キーの項目とそれ以外の項目の間には1対多の関係が常に存在する。図6の例では、支店コードと県コード、支店コードと地域コードの



間に1対多の関係がある。もちろんこれらの関係は、図3の処理303で説明した条件項目の選択時に、条件項目間の相互情報量を算出したときに、容易に自動的に発見することもできる。

【0159】これらの上位概念リストを用いてルールの生成時の木構造の探索を行うときには、追加すべき述語に対応する項目の上位概念または下位概念の項目がすでに条件節に含まれているかどうかをチェックする（直接の上位/下位概念だけでなく、上位の上位/下位の下位にあたる項目などもチェックする）。もし含まれていれば、その項目の述語は追加の対象とはしない。

【0160】次に、数値→記号値変換でファジィ分割を用いる場合の処理内容を説明する。

【0161】本実施例のルール生成処理では、図11のファジィ分割方法がもつ特徴を利用する。図11のファジィ分割では、1つの数値は高々2つの記号値に分配され、それらの和は必ず1.0になる（例えば、「小さい」0.3、「中くらい」0.7、「大きい」0.0のように）。また、図4のルール生成用ワークファイル405中の表現では、「小」を0に、「中」を1に、「大」を2に割りあてる。このとき、数値  $x$  は次の規則により1つの実数値  $y$  に変換される。

【0162】ファジィ記号値変換規則：もし「小さい」の割合が  $z$ 、「中くらい」の割合が  $(1-z)$  ならば、 $y = 1-z$

もし「中くらい」の割合が  $z$ 、「大きい」の割合が  $(1-z)$  ならば、

$$y = 2-z$$

ただし、 $0 \leq z \leq 1$  とする。

【0163】すなわち、0と1の間の実数値  $y$  は、「小」と「中」という記号の配分比を表すものと約束する。図11のファジィ分割の特徴から、1つの数値が3つ以上の記号値に分配されることはなく、また和が1.0になることから、この定義で矛盾は生じない。

【0164】上記の方法でファジィ的に記号化された事例を  $(y_1, y_2, \dots, y_n)$  とする。ルール生成処理において事例の数を計算するのは、ルールの評価尺度を計算するときである。クリスプな分割のときは、ある条件節を満たす事例がいくつあるかをカウントする処理になる。ファジィ記号化では次のようになる。

【0165】 $x_1$  が「小さい」に該当する事例は、 $0 \leq y_1 \leq 1$  の値をもつ事例で、その重みは  $1-y_1$  になる。

【0166】 $x_1$  が「中くらい」に該当する事例は、 $0 \leq y_1 \leq 2$  の値をもつ事例で、その重みは  $1-|1-y_1|$  になる。

【0167】 $x_1$  が「大きい」に該当する事例は、 $1 \leq y_1 \leq 2$  の値をもつ事例で、その重みは  $y_1-1$  になる。

【0168】もちろん、条件節述語が複数ある場合には、上記の重みをそれぞれの述語に関して掛け合わせたものが全体の事例の重みになる。ルール評価尺度の定義

に現れる  $N(A)$  や  $N(A \& B_i)$  などは、これらの事例重みの和として考え、実数値として取り扱うものとする。

【0169】上記の処理でファジィ分割により記号値に変換した項目をルール生成処理で使用すれば、分割における境界付近のデータを正しく扱うことができるので、生成されたルールを高精度化することができる。

【0170】図4のワークファイル405中のデータでは、値-1は欠損値を表すものとする。ルール生成処理には数値属性そのものは使用せず、また記号値は0から順番に割り当てて行くので、この値が混乱を起こすことはない。

【0171】欠損値自身をルール生成に使用するかどうかは、図3の処理303における欠損値対応処理か、またはユーザが指定した内容で決まるものとする。ルール生成に使用する場合には、「もし  $x$  が欠損値で、…」というルールが許される。欠損値をルール生成に使用しない場合には、単に欠損値は無視される（すなわちどの記号にも一致しない記号として扱われる）。欠損値は特別な値なので、前項のファジィ記号化の規則は適用されない。

【0172】本実施例では、ルール生成処理がどこまで進んだかを示す進行状況を表示する。ルール生成処理では、本来調べるべきルールの個数はあらかじめ算出できるから、その内の何パーセントを調べたかという数字を、適当なタイミングで表示するものとする。枝刈りなどで処理をスキップされたルールも当然カウントする。

【0173】ルール生成処理では、ユーザ割り込みを受け付けるものとする。上述のルール生成進行状況表示に、中断コマンドを受け付けるボタン1401を用意しておく。図14に、ルール生成処理進行状況および中断ボタン1401の表示例を示す。

【0174】ボタン1401がクリックされ処理の中断が指定された場合、次の動作は「続行」か「終了」かのどちらかになる。続行の場合にはそのまま処理を続け、終了の場合にはルール生成処理を終了し、現時点までに生成されたルールを表示するものとする。

【0175】図15に、以上説明したルール生成処理の処理フロー図を示す。候補ルール集合のクリアおよび  $\mu_{th}$  のゼロクリアなどの初期化（1501）の後、すべての候補についてステップ1503～1512を繰り返す。まず、評価尺度  $\mu$  を計算し、 $\mu$  と  $\mu_{th}$  とを比較する（1504）。 $\mu > \mu_{th}$  なら、候補ルール集合に当該候補ルールを追加し（1505）、 $\mu_{th}$  に候補ルール集合中で最小の評価尺度をセットする（1506）。 $\mu \leq \mu_{th}$  なら、サブノードに関する評価尺度の最大値  $\mu_{max}$  を計算し（1507）、 $\mu_{max}$  と  $\mu_{th}$  とを比較する。 $\mu_{max} < \mu_{th}$  なら、サブノードをすべて評価済みとする（1509）。

【0176】次に、ルール生成処理進行状況を表示し



(1510)、処理中断が指示されたかどうか判定する(1511)。中断が指示されていたら、ループから脱出する(1512)。すべての候補についてループ処理が終了するか、ステップ1512でループから脱出したときは、候補ルール集合を出力して(1513)、処理を終了する。

【0177】次に、図4のルール学習処理実行モジュール403で実行される処理について説明する。

【0178】図16は、学習すべきメンバシップ関数の形状パラメータを示す。本実施例では、学習によって調節されるパラメータは分割点 $\alpha_1$ と $\alpha_2$ とする。すなわち、ルール生成処理で生成されたルールに対して、最もルール評価尺度が大きくなる分割点 $\alpha_1$ と $\alpha_2$ を求める。ただし、ファジィ分割を用いたルール生成処理において、1つの数値は高々2つの記号値に分配され、それらの和は必ず1.0になるという性質を仮定しているため、その仮定を満足する範囲でのみ $\alpha_1$ と $\alpha_2$ の値を変化させるものとする。図16には、学習後のメンバシップ関数形状の例を示した。

【0179】具体的な学習は次のようなステップで行う。

ステップ1：学習すべきルールを指定する。どのルールを学習するかをユーザが指定する。

ステップ2：ステップ3～ステップ8を学習終了条件が成立するまで繰り返す。

ステップ3：各パラメータの変更量 $\Delta\alpha_i=0$  ( $i=1, 2, \dots, m$ ) とする。 $m$ は調節すべき分割点パラメータの総数である。

ステップ4：指定されたルールのそれぞれに対してステップ5～ステップ7を繰り返す。

【0180】ステップ5：処理中のルールの評価尺度 $\mu$ を算出すると同時に、 $\partial N(A)/\partial\alpha_i$ 、 $\partial N(B)/\partial\alpha_i$ 、 $\partial N(A\&B)/\partial\alpha_i$ を算出する。 $N(A)$ は条件部Aを満足する(ファジィ分割を考慮した)事例の数であり、 $\alpha_i$ を微小量 $\Delta\alpha_i$ 変化させたときの $N(A)$ の変化量 $\Delta N(A)$ は、図16のメンバシップ関数形状から容易に求めることができる。 $\partial N(B)/\partial\alpha_i$ 、 $\partial N(A\&B)/\partial\alpha_i$ の値も同様に求める。

ステップ6： $\mu$ の定義式(4)とステップ5で求めた値から、 $\partial\mu/\partial\alpha_i=\partial N(A)/\partial\alpha_i\cdot\partial\mu/\partial N(A)+\partial N(B)/\partial\alpha_i\cdot\partial\mu/\partial N(B)+\partial N(A\&B)/\partial\alpha_i\cdot\partial\mu/\partial N(A\&B)$ なる関係式を用いて、 $\partial\mu/\partial\alpha_i$  ( $i=1, \dots, m$ ) を求める。

【0181】ステップ7：全ての $i$ に対して、 $\Delta\alpha_i=\Delta\alpha_i+\partial\mu/\partial\alpha_i$ とする。

ステップ8：全ての $i$ に対して、 $\alpha_i=\alpha_i+\lambda\cdot\Delta\alpha_i$ とする。 $\lambda$ はあらかじめ定められた学習係数である。

ステップ9：学習済みの $\alpha_i$ の値を図4のルールファイル407に出力する。

【0182】ステップ2における終了条件としては、例

えば

- ・指定された回数ステップ3～ステップ8を繰り返した
- ・全ての $i$ に対して $|\Delta\alpha_i|<\epsilon$ が成立した( $\epsilon$ はあらかじめ定められた正の定数)

- ・ルール評価尺度 $\mu$ の変化 $|\Delta\mu|<\epsilon$ が成立した

などの条件、または上記条件の組み合わせ、などを用いることができる。

【0183】以上の処理により、1つまたは複数のルールの評価尺度を最大化するようなメンバシップ関数パラメータを学習により自動的に求めることができる。これにより、データ間に存在する規則性・因果関係をより精度よく抽出することが可能になるとともに、規則性を説明するために最適な項目値の分割点自身を最適に決定することができるようになるという効果がある。また、この学習によれば、生成ルール全体またはその部分集合の評価尺度の和を最適化することが可能なため、ユーザにとって意味のあるルールだけを学習によって最適化することができ、より利用価値の高いルールを生成することができるという効果もある。

【0184】

【発明の効果】以上説明したように、本発明によれば、次のような効果を実現することができる。

【0185】まず、データ間に存在する規則性・因果関係をユーザにとって理解しやすいルールという形で抽出できるため、従来の方法では困難であった大量データの有効活用が可能になるという効果がある。特に、製造プロセスにおける検査データの分析による不良発生要因の解析や、顧客データベースの分析による新商品のマーケティング戦略の立案などにおいて、有効な情報をデータベースから抽出することができるようになるという効果がある。

【0186】さらに、データベース中に含まれる項目のうち、ルール生成に使用する項目を自動的に選択でき、またその他の処理パラメータもデフォルト値が用意されているため、データの性質に関して詳細な知識を持たないユーザでも、分析に当たって複雑な前処理をすることなしに、データ分析処理をほぼ全自動的に実行することができ、きわめて使い勝手のよい分析手法を提供することが可能になるという効果がある。

【0187】また、データベースに含まれる数値データは、記号値に変換されてルールに使用されるため、雑音の大きなデータに対してもマクロな構造に対応した規則性を発見でき、それだけユーザにとって利用価値の高い結果を生成できるという効果がある。

【0188】上記記号値への変換をファジィ的に行うようにすれば、記号値へ変換することによる情報の損失を最小限に押さえ、精度のよいルールを抽出できるという効果がある。

【0189】また、データ中に欠損値が含まれる場合でも、その欠損値に意味があるとしてルールに使用するこ

【0194】また、記号値として与えられた項目の間の依存関係をデータ定義情報やデータ自身の分布から求め、ルール生成において冗長なルールを生成しないようにすることができるため、より単純でユーザの理解が容易なルールを、高速に生成することができるという効果がある。

【符号の説明】

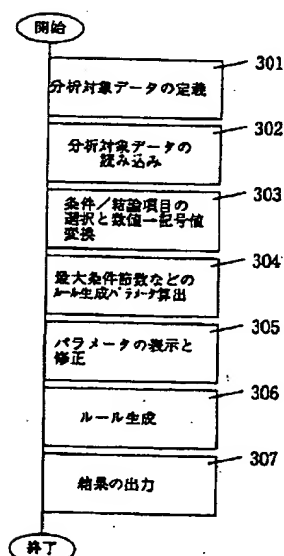
401…ユーザインタフェースモジュール、402…ルール生成処理実行モジュール、403…ルール学習処理実行モジュール、404…分析対象データベース、405…ルール生成用ワークファイル、406…ルール生成パラメータファイル、407…ルールファイル。

[illegible]

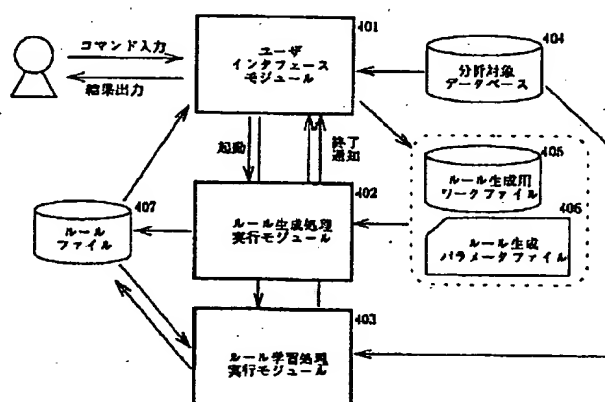
【図2】

顧客番号	氏名	年齢	性別	住所	支店コード	県コード	地域コード	預金残高	ローン残高	クレジットカード	...	商品購入
0158001	本多一郎	35	男	東京都世田谷区	108	25	3	250	1200	あり		あり
2048002	細川和夫	58	男	横浜市緑区	121	26	3	1200	510	あり		あり
3398024	朝倉知美	-	女	兵庫県芦屋市	257	37	5	120	340	なし		なし
1259842	上杉道彦	42	男	-	292	42	6	840	920	なし		あり
...	...	...	...	...	...	...	...	...	...	...	...	...
0912637	原川昌平	42	男	福岡県福岡市	359	46	7	780	0	あり		なし

【図3】



【図4】



【図5】

顧客番号	氏名	年齢	性別	住所	支店コード	県コード	ローン残高	預金残高	クレジットカード
0158001	本多一郎	35	男	東京都世田谷区	108	250	1200	あり	

(a) 顧客情報テーブル

支店コード	県コード	地域コード	支店長氏名
108	25	3	松平明夫

(b) 支店情報テーブル

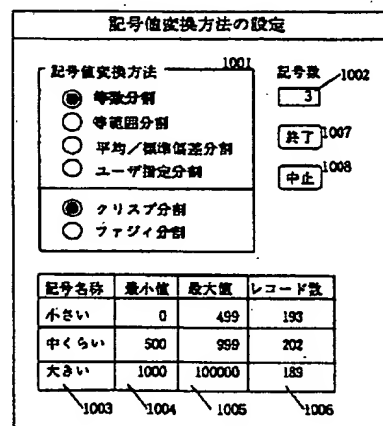
顧客番号	商品購入
0158001	あり

(c) 商品購入履歴テーブル

【図6】

顧客番号	氏名	年齢	性別	住所	支店コード	預金残高	ローン残高	クレジットカード
0158001	本多一郎	35	男	東京都世田谷区	108	250	1200	あり

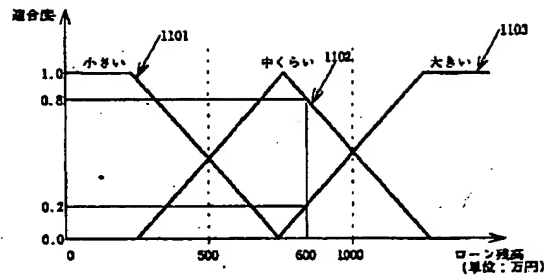
【図10】



【図7】

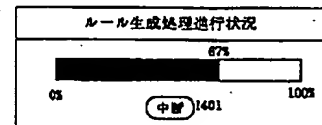
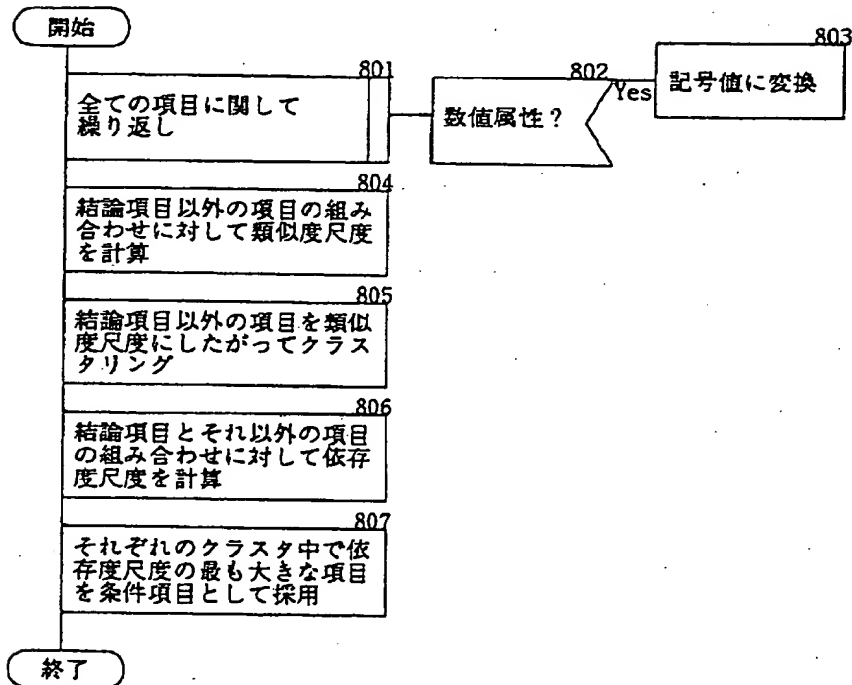
条件項目・結論項目の指定			
未使用項目一覧 顧客番号 氏名 年齢 性別 住所 支店コード 県コード 地域コード 預金残高 ローン残高 クレジットカード 商品購入	自動選択	704	条件項目 702
	すべて追加 >	705	
	追加 >	706	
	< 削除	707	
	<< すべて削除	708	結論項目 703
	追加 >	709	
	< 削除	710	
			終了 711
			中止 712

【図11】



【図14】

【図8】



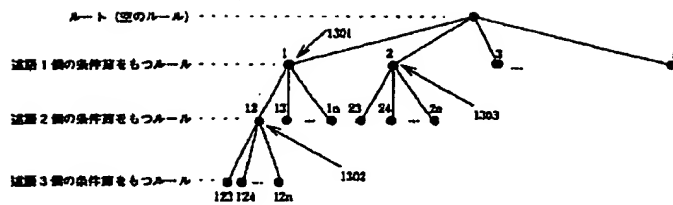
【図9】

記号値変換方法一覧			
項目名	記号型	記号名称	変換方法
年齢	5	20未満   20代   30代   40代   50以上	ユーザ指定(7桁)
預金残高	3	少ない   ふつう   多い	等数分割(7桁)
ローン残高	3	小さい   中くらい   大きい	等値区分割(7桁)

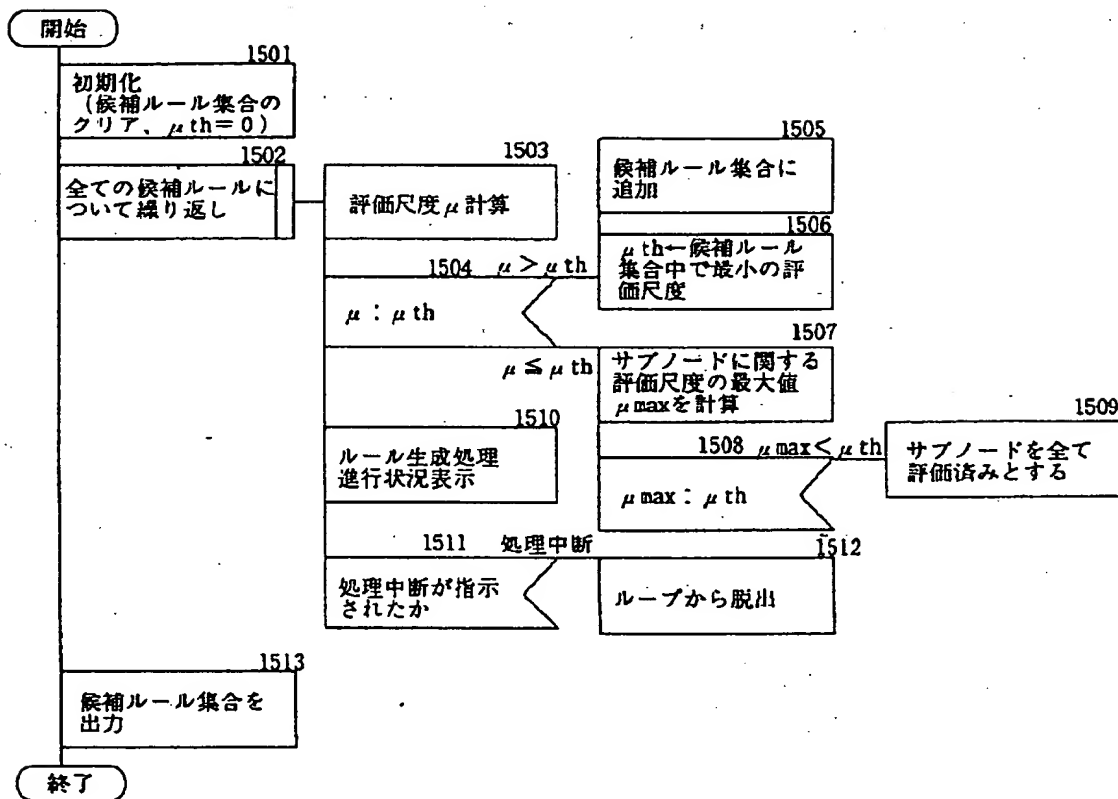
【図12】

ルール生成パラメータの指定	
生成ルール数	24 1201
全体カバー率	1202
最大条件数	3 1203
最小カバー率(%)	4 1204
カバー率優先係数	1.0 1205
データファイル名: 顧客情報	1206
レコード数: 2354	1207
条件項目数: 8	1208
補償部: 商品購入: あり	

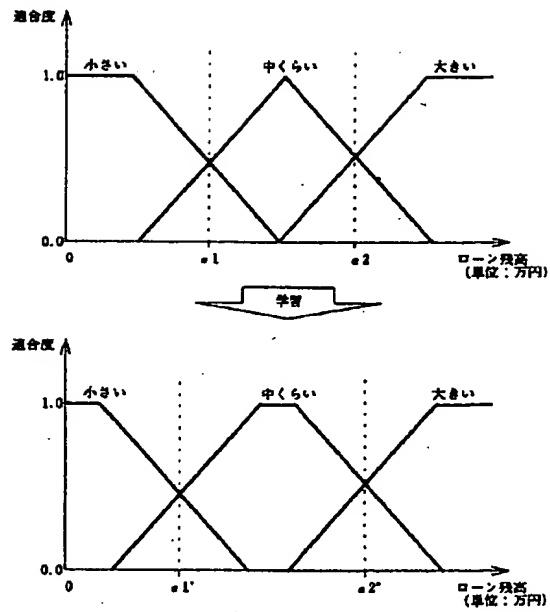
【図 13】



【図 15】



【図16】



フロントページの続き

(72)発明者 伊藤 幸康  
神奈川県横浜市戸塚区戸塚町5030番地 株  
式会社日立製作所ソフトウェア開発本部内

(72)発明者 高橋 ヨリ  
神奈川県横浜市戸塚区戸塚町5030番地 株  
式会社日立製作所ソフトウェア開発本部内